



THE UNIVERSITY *of* EDINBURGH

## Edinburgh Research Explorer

### Meta-Analysis of 1,200 Transcriptomic Profiles Identifies a Prognostic Model for Pancreatic Ductal Adenocarcinoma

**Citation for published version:**

Sandhu, V, Latori, KJ, Borgida, A, Lungu, I, Bartlett, J, Hafezi-bakhtiari, S, Denroche, RE, Jang, GH, Pasternack, D, Mbaabali, F, Watson, M, Wilson, J, Kure, EH, Gallinger, S & Haibe-kains, B 2019, 'Meta-Analysis of 1,200 Transcriptomic Profiles Identifies a Prognostic Model for Pancreatic Ductal Adenocarcinoma', *JCO Clinical Cancer Informatics*, no. 3, pp. 1-16. <https://doi.org/10.1200/CCI.18.00102>, <https://doi.org/10.1200/CCI.18.00102> JCO Clinical Cancer Informatics -

**Digital Object Identifier (DOI):**

[10.1200/CCI.18.00102](https://doi.org/10.1200/CCI.18.00102)

[10.1200/CCI.18.00102](https://doi.org/10.1200/CCI.18.00102) JCO Clinical Cancer Informatics -

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

JCO Clinical Cancer Informatics

**Publisher Rights Statement:**

This is a pre-copyedited, author-produced version of an article accepted for publication in The Journal of Clinical Oncology following peer review. The version of record "Meta-Analysis of 1,200 Transcriptomic Profiles Identifies a Prognostic Model for Pancreatic Ductal Adenocarcinoma " is available online at: <https://ascopubs.org/doi/pdf/10.1200/CCI.18.00102>

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# **Meta-analysis of 1,200 transcriptomic profiles identifies a prognostic model for pancreatic ductal adenocarcinoma**

## **Authors**

Vandana Sandhu<sup>1,2</sup>, Knut Jorgen Labori<sup>3</sup>, Ayelet Borgida<sup>4</sup>, Ilinca Lungu<sup>5</sup>, John Bartlett<sup>5</sup>, Sara Hafezi-Bakhtiari<sup>6</sup>, Rob Denroche<sup>8</sup>, Gun Ho Jang<sup>8</sup>, Danielle Pasternack<sup>7</sup>, Faridah Mbaabali<sup>7</sup>, Matthew Watson<sup>7</sup>, Julie Wilson<sup>8</sup>, Elin H. Kure<sup>2,9</sup>, Steven Gallinger<sup>8,10,11</sup>, Benjamin Haibe-Kains<sup>1,8,12,13</sup>

## **Affiliations**

1. Princess Margaret Cancer Centre, University Health Network, Toronto, M5G 1L7, Canada
2. Department of Cancer Genetics, Institute for Cancer Research, Oslo University Hospital - The Norwegian Radium Hospital, Oslo, Norway
3. Department of Hepato-Pancreato- Biliary Surgery, Oslo University Hospital, Oslo, Norway
4. Zane Cohen Centre for Digestive Diseases, Mount Sinai Hospital, Toronto
5. Diagnostic Development, Ontario Institute of Cancer Research
6. Department of Pathology, University Health Network, Toronto, M5G 1L7, Canada
7. Department of Genomics, Ontario Institute of Cancer Research
8. PanCuRx Translational Research Initiative, Ontario Institute of Cancer Research
9. Department of Natural Sciences and Environmental Health, Faculty of Technology, Natural Sciences and Maritime Sciences, University of Southeastern Norway, Bø in Telemark, Norway
10. Wallace McCain Centre for Pancreatic Cancer, University Health Network, Toronto, M5G 1L7, Canada
11. Hepatobiliary/pancreatic Surgical Oncology Program, University Health Network, Toronto, M5G 1L7, Canada
12. Department of Medical Biophysics, University of Toronto, Toronto, M5G 1L7, Canada
13. Department of Computer Science, University of Toronto, Toronto, M5S 2E4, Canada

**\*Corresponding author:** Benjamin Haibe-Kains

101 College Street, PMCRT 11-310, M5G1L7, Toronto, Canada

[benjamin.haibe.kains@utoronto.ca](mailto:benjamin.haibe.kains@utoronto.ca), Phone: +1 (416) 581-7628

**Running head title:** Meta analysis to predict overall survival for pancreatic cancer

## **ABSTRACT**

### **Purpose**

With a dismal 8% median 5-year overall survival (OS), pancreatic ductal adenocarcinoma (PDAC) is highly lethal. Only 10-20% of patients, which are eligible for surgery, and over 50% of these will die within a year of surgery. Building a molecular predictor of early death would enable the selection of PDAC patients at high risk.

### **Materials and Methods**

We developed the Pancreatic Cancer Overall Survival Predictor (PCOSP), a prognostic model built from a unique set of 89 PDAC tumors where gene expression was profiled using both microarray and sequencing platforms. We used a meta-analysis framework based on the binary gene pair method to create gene expression barcodes robust to biases arising from heterogeneous profiling platforms and batch effects. Leveraging the largest compendium of PDAC transcriptomic datasets to date, we show that PCOSP is a robust single-sample predictor of early death ( $\leq 1$  yr) after surgery in a subset of 823 samples with available transcriptomics and survival data.

### **Results**

The PCOSP model was strongly and significantly prognostic with a meta-estimate of the area under the receiver operating curve (AUROC) of 0.70 ( $P=2.6e-22$ ) and hazard ratio (HR) of 1.95(1.6-2.3) ( $P=1.4e-04$ ) for binary and survival predictions, respectively. The prognostic value of PCOSP was independent of clinicopathological

parameters and molecular subtypes. Over-representation analysis of the PCOSP 2619 gene-pairs (1070 unique genes) unveiled pathways associated with Hedgehog signalling, epithelial mesenchymal transition (EMT) and extracellular matrix (ECM) signalling.

### **Conclusion**

PCOSP could improve treatment decision by identifying patients who will not benefit from standard surgery/chemotherapy and may benefit from a neoadjuvant approach.

### **INTRODUCTION**

Pancreatic ductal adenocarcinoma (PDAC) is a highly lethal malignancy with 5-year overall survival rate less than 8%<sup>1</sup>. The majority of patients (> 80%) are inoperable due to locally advanced or metastatic disease at time of diagnosis. While surgical resection is the key to curative treatment, it rarely results in long-term survival<sup>2</sup>. Hence, completion of multimodality treatment - surgery combined with adjuvant or neoadjuvant chemotherapy- is the standard of care for treatment of PDAC. However, even after surgical resection with curative intent, median survival does not exceed 28 months and half of those who undergo surgery develop recurrent disease, and die within a year after surgery <sup>2-4</sup>. Therefore, there is a need for a robust prognostic model to identify patients with high risk of early death based on molecular profiles of their tumors. Such a prognostic model would assist clinicians in identifying patients who might not benefit from surgery and standard adjuvant chemotherapy and may benefit from a neoadjuvant approach. Neoadjuvant treatment sequencing is the only alternative strategy and may guide selection of patients for surgery and help to

identify those patients with progressive disease for whom an operation has little oncologic benefit.

Various clinical factors are prognostic following PDAC surgery such as lymph node metastasis status<sup>5</sup>, tumor grade<sup>6</sup>, margins<sup>7</sup>, degree of differentiation<sup>8</sup> and protein biomarker CA-19-9<sup>9</sup>. However, the prognostic value of these clinical variables are insufficient to accurately stratify patients based on risk of disease recurrence<sup>10, 11</sup>. With the advent of high-throughput next-generation molecular profiling technologies, multiple studies have released transcriptomic profiles of PDAC to the public domain. These gene expression profiles have been leveraged to identify molecular subtypes of PDACs<sup>12-16</sup>. While overlap between these subtypes<sup>15</sup> supports the biological relevance of these published classification schemes<sup>15</sup>, they have not been designed to optimize prognostic value.

Previously published prognostic models were developed from small number of samples lacking proper validation in multiple datasets<sup>17-21</sup>. Attempts have been made recently to build a prognostic gene signature using pooled samples from multiple cohorts to identify patients at high risk of short-term survival post surgery<sup>22-24</sup>. However, they used samples profiled using either array or sequencing based method as the learning cohort, therefore the classifiers may perform better for subjects whose samples were profiled using only one of the two platforms.

To address these issues, we took advantage of a unique set of 89 PDACs profiled using both microarray and sequencing technologies to develop the Pancreatic Cancer Overall Survival Predictor (PCOSP) model. Using an independent set of PDAC transcriptomic profiles from 823 primary resected patients, we show

that PCOSP is a robust single-sample predictor of early death ( $\leq 1$  yr) after surgery, which could be used as as a potential tool to assist clinicians in decision making. with a meta-estimate of the area under the receiver operating characteristics curve of 0.70 ( $p=1.9e-18$ ). We also show that PCOSP is significantly prognostic (meta-estimate of hazard ratio of 1.95;  $p=2.6e-16$ ). Furthermore, we show that PCOSP performs significantly better than published prognostic models across microarray and sequencing datasets (Superiority test,  $P < 0.01$ ). Our results support PCOSP as a potential tool to assist clinicians in decision making.

## **MATERIALS AND METHODS**

The meta-analysis pipeline used to develop the PCOSP model and evaluate its prognostic value is provided in Figure 1.

### **Datasets**

We surveyed the literature and curated 17 datasets including 1,236 PDAC patients from public domain for which transcriptome data of PDAC are available (Supplementary Table S1). We further filtered samples based on the availability of overall survival (OS) and sample size ( $>10$ ) after dichotomization into high and low survival groups based on an OS cut-off of 1-year (Figure 2). This resulted in a total of four sequencing studies and seven array-based studies providing transcriptomic and clinical data for 1,001 PDAC patients. A total of 12,430 protein-coding genes commonly assessed across all the cohorts were used for further analysis. The different cohorts had similar clinical presentation, and were treated with curative

surgery followed by adjuvant chemotherapy, where 2/3rd of the patients completed multimodal treatment (i.e., surgery and adjuvant chemotherapy) (Supplementary Table S2).

### **Prognostic model**

To develop a robust predictor for early death, we used the gene expression profiles of 89 PDAC patient samples whose tumors have been profiled using both microarray and sequencing platforms within the ICGC cohort. Human research ethical approval were as mentioned in <sup>14</sup>. Approximately half of the patients of the training cohort which were eligible for surgery relapsed within 1 year, we used this threshold to predict PDAC patients with high risk of early death ( $\leq 1$  yr) post surgery. We excluded 7 samples from the training cohort as these patients were censored before one year of follow-up.

To make gene expression profiles comparable between the training and validation sets, we transformed the original gene expression profiles into binary gene pair barcodes. The advantages of considering pairs of genes with a binary value (“1” if expression of gene  $i >$  gene  $j$ , “0” otherwise) are; (i) it transforms the feature space in a way that mitigates platform biases and potential batch effects; (ii) it makes the model robust to any data processing that preserves the gene order<sup>25, 26</sup>. We implemented k-Top Scoring disjoint Pairs (k-TSP) classifier predictor<sup>27</sup> using the Wilcoxon rank sum method as filtering function in the *SwitchBox* package (version 1.12.0)<sup>28</sup>.

The decision rules are based on the relative ordering of gene expression values within the same sample, where the  $k$  top scoring gene pairs are used to build the classifier. The samples were resampled 1000 times, where 40 samples from each group were selected in each run to build a  $k$ -TSP model and the model was further tested on the 49 out-of-bag samples. The models were selected if the balanced accuracy was above 0.6 else the model was rejected. We then froze the parameters of the predictive model and validated it in the remaining compendium of independent datasets. The class probability of the sample was calculated as the frequency of sample predicted as one class divided by the total number of models.

### **Random classifiers**

To test whether the prognostic value of the PCOSP model could be achieved by random chance alone, we implemented two permutation tests. To test whether the gene expression profiles were associated with survival, we shuffled the actual class labels while maintaining the expression values. To test whether the gene pairs selected in the PCOSP model were robustly associated with survival, we randomly assigned genes to the  $k$ -TSP model and assessed its prognostic value. Both procedures were performed 1000 times. As a pre-validation set we compared the balanced accuracy of all the 1000 random models generated using both the approaches to PCOSP using the Wilcoxon rank sum test. Further, we trained the  $k$ -TSP classifier models from both approaches in the same way as we built our consensus PCOSP model. We then froze the parameters of the prognostic model



and validated it in the compendium of independent datasets, and compared the meta-estimates for both the models against the PCOSP model.

### **Early death prediction**

The meta analysis was performed for the PDAC sequencing cohorts, PDAC array-based cohorts and the overall combined cohorts to assess and statistically compare the performance of the PCOSP. The patient samples were dichotomized into two groups based on the outcome variable (time from surgery to death  $\leq 1$  year). Samples censored before 1 year of follow-up were excluded from the analysis of meta-estimate of the area under the receiver operating characteristics curve (AUROC). The AUROC plots the sensitivity vs. 1-Specificity and is used as a criterion to measure the discriminatory ability of the model<sup>29</sup>. The AUROC was computed using *pROC* package (version 1.10.0), and the p-value was estimated using the Mann-Whitney test statistics estimating whether the AUROC curve estimate is significantly different from 0.5 (random classifier). The meta-estimate of AUROC was estimated using the random effect model<sup>30</sup> implemented in *survcomp* package (version 1.26.0)<sup>31, 32</sup>.

### **Survival prediction**

Prognostic value and statistical significance of survival difference between the predicted classes were assessed using the D-Index, which is a robust estimate of the traditional Cox's hazard ratio, more precisely an estimate of the log hazard ratio comparing two equal-sized prognostic groups<sup>33</sup> and is a natural measure of separation

between two independent survival distributions under the proportional hazards assumption<sup>33</sup>. In addition, we used the concordance index (C-index) which estimates the probability that, for a random pair of patients, the PCOSP score for the patient with shorter survival is higher than the patient with longer survival<sup>34</sup>. Both the robust hazard ratio (HR) and the C-index were calculated using the *survcomp* package. The meta estimate of HR and C-index were calculated for the PDAC sequencing cohorts, the PDAC array-based cohorts and the combined PDAC sequencing and array-based cohorts using the random effect model<sup>30</sup> implemented in *survcomp* package. The patients were stratified into low- and high-risk group using median PCOSP score as a threshold. Kaplan Meier curves were plotted using *survminer* package (version 0.4.3)<sup>35</sup> in R and reported the P values from log-rank test.

### **Subtyping of PDAC cohorts**

The PDAC cohorts were classified into basal and classical transcriptomic subtype using the Moffitt classifier<sup>13</sup>.

### **Clinicopathological features based model to predict early death**

The clinical model was built by fitting the logistic regression model using common clinicopathological features i.e., age, gender, TNM status and tumor grade available from PCSI, ICGC-sequencing, ICGC-array, TCGA and OUH cohorts.

### **Gene set enrichment analysis**

To categorize genes in the PCOSP, we performed gene set enrichment analysis using RunGSAhyper function implemented in piano package (version 1.16.4)<sup>36</sup>. The genes selected in the PCOSP model (n=1,070) were compared against Gene Ontology (GO) gene sets, canonical pathways and hallmark gene sets in MSigDb<sup>37, 38</sup>, using as background the protein-coding genes commonly assessed across the gene expression profiling platforms in our data compendium. Enrichment p-values were corrected for multiple testing using the false discovery rate approach (FDR < 5%)<sup>39</sup>.

### **Comparison to existing classifiers**

We calculated the Birnbaum signature scores<sup>22</sup> and Chen signature scores<sup>23</sup> using the published coefficients of the 25 and 15 classifier genes, respectively, as weight parameter in the *sig.score* function implemented in the *genefu* R package (version 2.10.0)<sup>40</sup>. The Haider signature scores were used as courtesy of the author<sup>24</sup>. The C-index and HR were computed for the three classifiers using eight validation cohorts excluding the cohorts used for training by PCOSP and other classifiers in comparison. Further, we compared the meta-estimates of C-index of each classifier with PCOSP at P<0.05 (one-sided t-test) as implemented in *survcomp* package.

### **Research reproducibility**

Our code and documentation are open-source and publicly available through the PDACSurv GitHub repository ([github.com/bhklab/PDACSsurv](https://github.com/bhklab/PDACSsurv)). A detailed tutorial describing how to run our pipeline and reproduce our analysis results is available in

the GitHub repository. A virtual machine reproducing the full software environment is available on [Code Ocean](#). Our study complies with the guidelines outlined in <sup>41–43</sup>. All the data are available in the form of R package [MetaGxPancreas](#).

## RESULTS

### Overall survival predictive model

To predict the patients with early death ( $\leq 1$  year after surgery), the PCOSP model was trained on the 89 ICGC cohort samples profiled using both microarray and sequencing transcriptomic profiles (Supplementary Table S1). To develop a predictor that can be applied to multiple profiling platforms, we transformed the gene expression profiles into binary gene pairs ( $x=1$  if expression of gene  $i >$  gene  $j$ ,  $x=0$  otherwise) and used these transcriptomic barcodes in an ensemble of 1000  $k$ -TSP predictive models. The PCOSP score is subsequently calculated using the majority voting rule. We tested the prognostic value of PCOSP score in three independent sequencing cohorts, including the Pancreatic Cancer Sequencing Initiative (PCSI)<sup>44</sup>, TCGA-PAAD<sup>15</sup> and Kirby<sup>45</sup> cohorts, and seven independent array-based cohorts composed of ICGC-array (excluding the 89 samples used for training)<sup>46</sup>, UNC<sup>13</sup>, OUH<sup>47</sup>, Chen<sup>23</sup>, Zhang<sup>48</sup>, Winter<sup>49</sup> and Collisson cohorts<sup>12</sup> (Supplementary Table S1). We first tested the predictive value of early death by calculating the AUROC for each dataset separately (Figure 3A). PCOSP was significant overall (AUROC=0.70;  $P<2.6E-22$ ; Figure 3A), although higher in the datasets generated using sequencing platforms compared to microarrays (AUROC 0.72 vs 0.68 for sequencing and array datasets, respectively) at ( $P=0.09$ ) suggesting that RNA-sequencing might be a

better assay for PCOSP than microarray platforms. PCOSP was significantly predictive of early death in all cohorts ( $\text{AUROC} \in [0.67, 0.76]$ ;  $P < 0.05$ ) except the Winter and OUH cohorts ( $P > 0.48$ ) and was almost significant for the Collisson cohort ( $\text{AUROC} = 0.69$ ;  $P = 0.051$ ). To determine whether the early death predictive value of the PCOSP model can be achieved by random chance alone, we first computed meta-estimates of AUROC by randomly shuffling the class labels (early deaths) 1000 times and applying the same training procedure used for the PCOSP model. We observed that the gene expression profiles were significantly associated with survival as none of the random models could yield a predictive value greater or equal to PCOSP ( $p < 0.001$ ; Supplementary Figure S1A). We further tested whether the gene pairs selected in the PCOSP model were robustly associated with early death events, by randomly assigning genes to the PCOSP model. Again, we observed that the genes selected in PCOSP yielded significantly more predictive information than the models comprised of random genes ( $p < 0.001$ ; Supplementary Figure S1B), supporting the biological relevance of the PCOSP gene set.

### **Prognostic relevance of the PCOSP model**

To assess the prognostic value of the PCOSP model, we calculated the C-indices and HR using the overall survival data for all the cohorts. The C-index is significant overall ( $\text{C-index} = 0.63$ ,  $P = 1.8\text{E-}12$ ; Figure 3B). In agreement with the results of early death prediction, the PCOSP prognostic value was higher for the sequencing datasets when compared to the arrays arrays ( $\text{C-index} = 0.65$  ( $P < 3.8\text{E-}14$ ) vs  $0.61$  ( $P < 1.6\text{E-}12$ ) for sequencing and array datasets, respectively; Figure 3B). Similar to

the C-index, the PCOSP HR was strong and significant overall (HR =1.95,  $P=1.4E-04$ ; Figure 3C), and stronger for the sequencing datasets (HR = 2.24 vs 1.83; Figure 3C). To assess whether the prognostic value of PCOSP depends on PDAC molecular subtypes, we stratified PDAC samples into the basal and classical subtypes using Moffitt classifier and calculated meta-estimates of C-index and HR (Supplementary Figures S2A and S2B). We found that PCOSP was prognostic in validation cohorts independently of molecular subtypes. We further tested whether PCOSP prognostic value was complementary to clinicopathological parameters and molecular subtypes by fitting both a multivariate Cox proportional hazard model to predict survival and a logistic regression model to predict binary outcome (death >1yr or not) (Supplementary Table S3).

To further illustrate the prognostic value of PCOSP, we stratified the patients into low- and high-risk group and plotted the KM curves for each cohort (Figure 4A-4J). The OS were significantly different between the risk groups for all the sequencing cohorts and 2 microarray cohorts ( $P<0.05$ ) and borderline significant for 3 microarray cohorts ( $0.05\leq P<0.10$ ; Figure 4A-4J); with 10-month difference in median OS between risk groups.

### **Clinicopathological model to predict overall survival**

The logistic regression model fitted using these clinicopathological features was used to predict early death of PDAC patients. The clinicopathological model was not significant overall (C-index=0.55;  $P=0.17$ ; Figure 5A). Contrary to PCOSP, the

clinicopathological model was not predictive in the sequencing cohort (C-index=0.53 and 0.58 with  $P=0.75$  and  $0.05$  for the sequencing and the array datasets, respectively; Figure 5A). Only nodal status, tumor grade and molecular classes were significant in the univariate analysis (Supplementary Table S3). We compared the prognostic value of the clinicopathological model against PCOSP (Figure 5B,C). PCOSP was significantly more prognostic than the clinicopathological model (one-sided t-test  $P < 0.01$ ; Figure 5D).

### **Comparison with published prognostic models**

We compared the prognostic value of PCOSP to three published PDAC prognostic models, referred to as Birnbaum<sup>22</sup>, Chen<sup>23</sup> and Haider<sup>24</sup>. The overall prognostic value of the three published models was significant (Figure 6A,C). PCOSP significantly outperformed published prognostic models in all cases ( $P < 0.05$ , Figure 6C,D); except for the HR of the Chen classifier where the superiority of the PCOSP prognostic value showed a trend to significance (one sided t-test  $P=0.10$ ).

### **Pathway analysis of prognostic genes**

Gene enrichment analysis for PCOSP signature genes ( $n=1,070$ ) was performed using hypergeometric test using the hallmarks gene sets, GO molecular function, GO cellular component terms and canonical pathways in MSigDb<sup>37</sup>. The Extracellular matrix (ECM), Epithelial Mesenchymal transition (EMT) and hedgehog signalling pathway genes were enriched in the PCOSP model at false discovery rate (FDR)

<5%. The complete list of GO terms and pathways significantly enriched in the PCOSP model are listed in Supplementary Table S4A- 4D.

## **DISCUSSION**

We performed a meta-analysis of the transcriptomic profiles of 1,236 PDAC patients and developed PCOSP, a new prognostic model to identify patients with high risk of early death after surgery. The model is built from a unique set of 89 patients profiled using both array-based and sequencing platforms, and validated on a compendium of ten independent datasets, including 823 patients. The prognostic value of the PCOSP model was highly significant for both early death ( $\leq 1$  year) and overall survival ( $P < 0.001$ ; Figure 3).

Contrary to published prognostic signatures fitted on small number of samples and lacking validation in large independent datasets<sup>17-21</sup>, PCOSP has been trained and validated on a large compendium of datasets. Comparison of PCOSP with existing classifiers<sup>22-24</sup> showed that the Birnbaum, Chen and Haider models yielded significant but significantly weaker prognostic value than PCOSP (Figure 6C,D). Importantly, PCOSP performs significantly better than existing classifiers for both microarray and sequencing platforms, likely due to simplifying the continuous expression space into binary pair barcodes. This enables PCOSP to be used as a single sample predictor robust to profiling platforms, potential batch effects and normalization methods compared to other classifiers.



Comparison of PCOSP against known prognostic clinicopathological variables showed that PCOSP outperformed the clinicopathological model in predicting early death (Figure 5). PCOSP prognostic value was significant, even after adjusting for molecular subtyping (classical vs basal) and clinicopathological parameters (age, sex, TNM status, differentiation grade of tumor and molecular classes) (Supplementary Figure S2A,B and Supplementary Table S3).

The PCOSP model incorporates 2,619 unique gene pairs, totalling 1,070 unique genes. Functional analysis of 1,070 genes showed enrichment of Hedgehog signalling, ECM and EMT pathway. Numerous studies have suggested the involvement of EMT in invasion and metastasis of PDAC<sup>50</sup>. EMT enhances cell motility through loss of cell-cell adhesion, escaping from extracellular matrix and overcoming the apoptosis process<sup>50</sup>. The ECM and EMT pathways are not only associated with the metastatic spread of tumor but also with chemoresistance that leads to worse survival<sup>51</sup>.

PDAC is a heterogeneous and genetically highly complex disease, supporting the molecular<sup>13, 14</sup> and morphological<sup>52</sup> characterization of a given tumor as an important cornerstone for the development of future therapies. We provide the largest compendium of 17 PDAC datasets as a gold standard for future PDAC analyses. The new meta-analysis framework implemented in PCOSP maximizes robustness and performance across the cohorts. In order to implement PCOSP as a clinical assay, we tested different feature set sizes for the k-TSP models and compared the performance of the reduced models. We achieved accuracy comparable to the 1,070 gene-PCOSP model by including only 256 unique genes,

supporting the potential of a smaller PCOSP based useful in the clinic (Supplementary Figure S3). Endoscopic ultrasound (EUS) biopsies could be utilized prior to curative surgery to estimate the prognosis of PDAC patients using PCOSP. This may assist clinicians in the selection of patients for surgery and help to identify those patients with high risk progressive disease for whom an operation has little oncologic benefit.

The current study has potential limitations. First, there are inherent tumor sample collection biases as the different datasets were collected and sampled at different centers. The levels of tumor cellularity varied highly across cohorts as PCSI and Collisson datasets were generated using laser microdissection prior to sequencing, Kirby and Chen datasets were macrodissected, while TCGA, ICGC, OUH, Zhang and Winter datasets used bulk tumors for profiling. Second, the transcriptomic profiles in our data compendium were generated using different gene expression profiling technologies for sequencing (Illumina HiSeq 2000/2500) and microarray platforms (Agilent, Affymetrix, and Illumina). Third, all samples were normalized using the published processing methods, which depend on the profiling platforms (Supplementary Table S2). Despite these limitations, PCOSP yielded robust prognostic value across the heterogeneous datasets, indicating that the gene expression barcode transformation is robust to the inevitable biases present in large meta-analyses. However, exploring other factors like germline variants, epigenetics, copy number alterations, non-coding RNAs, protein abundance as well as epidemiological and environmental factors will be necessary to further improve the prediction accuracy of predictive models.

The lack of available clinical and treatment information across the cohorts is also a limiting factor in our meta-analysis. However, comparison of cohort specific clinical information for the cohort were not significantly different across the cohorts (Supplementary table S2). During the time period of sample collection, standard of care treatment for PDAC was curative-intent surgery followed by adjuvant chemotherapy with gemcitabine or 5-FU. New approaches using doublet and triplet chemotherapy regimens are now standard of care in the palliative setting and randomised trials using these agents in the adjuvant setting will be reported shortly. Neoadjuvant therapy is also being evaluated in many centres. Thus, heterogeneity in treatment is expected within and between different cohorts, we will need to test our PCOSP model using new clinical datasets, or preferably within the context of randomized trials.

## **CONCLUSION**

We leveraged the largest compendium of PDAC transcriptomes to develop PCOSP, a prognostic model identifying PDAC patients at high risk of early death independently of, and superior to, clinicopathological features and molecular subtypes. PCOSP may be useful in the clinical setting as a single sample classifier to identify patients who could be at higher risk of early death following surgery and adjuvant chemotherapy, potentially facilitating treatment decisions, including the use of neoadjuvant chemotherapy as an alternative treatment strategy for these patients.

**List of Abbreviations:**

AUROC: Area under the receiver operating curve, GO: Gene annotation, OS: Overall survival, PCOSP: Pancreatic cancer overall survival predictor, PDAC: Pancreatic ductal adenocarcinoma, TSP: Top scoring pairs.

**Declarations:****Open Access**

This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

**Competing interests**

The authors declare that they have no competing interests.

**Acknowledgements**

This study was conducted with the support of the Ontario Institute for Cancer Research (OICR, PanCuRx Translational Research Initiative) through funding

provided by the Government of Ontario (Ministry of Research, Innovation, and Science), and a charitable donation from the Canadian Friends of the Hebrew University (Alex U. Soyka). V.S was supported by grants from The Radium Hospital Foundation, Oslo University Hospital, and the PanCuRx Translational Research Initiative at the OICR. B.H.K was supported by the Gattuso Slaight Personalized Cancer Medicine Fund at Princess Margaret Cancer Centre, the Canadian Institutes of Health Research, the Natural Sciences and Engineering Research Council of Canada, and the Ministry of Economic Development and Innovation/Ministry of Research & Innovation of Ontario (Canada). We thank Dr. Syed Haider for courteously providing the prediction scores from his classifier for comparison with PCOSP. We thank all the patients who participated in the study.

## REFERENCES

1. Siegel RL, Miller KD, Jemal A: Cancer statistics, 2017. *CA Cancer J Clin* 67:7–30, 2017
2. Winter JM, Brennan MF, Tang LH, et al: Survival after resection of pancreatic adenocarcinoma: results from a single institution over three decades. *Ann Surg Oncol* 19:169–175, 2012
3. Labori KJ, Katz MH, Tzeng CW, et al: Impact of early disease progression and surgical complications on adjuvant chemotherapy completion rates and survival in patients undergoing the surgery first approach for resectable pancreatic ductal adenocarcinoma – A population-based cohort study. *Acta Oncol* 55:265–277, 2016
4. Neoptolemos JP, Palmer DH, Ghaneh P, et al: Comparison of adjuvant gemcitabine and capecitabine with gemcitabine monotherapy in patients with resected pancreatic cancer (ESPAC-4): a multicentre, open-label, randomised, phase 3 trial. *Lancet* 389:1011–1024, 2017
5. Slidell MB, Chang DC, Cameron JL, et al: Impact of total lymph node count and lymph node ratio on staging and survival after pancreatectomy for pancreatic adenocarcinoma: a large, population-based analysis. *Ann Surg Oncol* 15:165–174, 2008
6. Lüttges J, Schemm S, Vogel I, et al: The grade of pancreatic ductal carcinoma is an independent prognostic factor and is superior to the immunohistochemical

assessment of proliferation. *J Pathol* 191:154–161, 2000

**7.** Richter A, Niedergethmann M, Sturm JW, et al: Long-term results of partial pancreaticoduodenectomy for ductal adenocarcinoma of the pancreatic head: 25-year experience. *World J Surg* 27:324–329, 2003

**8.** Imaoka H, Shimizu Y, Mizuno N, et al: Clinical characteristics of adenosquamous carcinoma of the pancreas: a matched case-control study. *Pancreas* 43:287–290, 2014

**9.** Tas F, Karabulut S, Ciftci R, et al: Serum levels of LDH, CEA, and CA19-9 have prognostic roles on survival in patients with metastatic pancreatic cancer receiving gemcitabine-based chemotherapy. *Cancer Chemother Pharmacol* 73:1163–1171, 2014

**10.** Le N, Sund M, Vinci A, et al: Prognostic and predictive markers in pancreatic adenocarcinoma. *Dig Liver Dis* 48:223–230, 2016

**11.** Martinez-Useros J, Garcia-Foncillas J: Can Molecular Biomarkers Change the Paradigm of Pancreatic Cancer Prognosis? *Biomed Res Int* 2016:4873089, 2016

**12.** Collisson EA, Sadanandam A, Olson P, et al: Subtypes of pancreatic ductal adenocarcinoma and their differing responses to therapy. *Nat Med* 17:500–503, 2011

**13.** Moffitt RA, Marayati R, Flate EL, et al: Virtual microdissection identifies distinct tumor- and stroma-specific subtypes of pancreatic ductal adenocarcinoma. *Nat*

Genet 47:1168–1178, 2015

- 14.** Bailey P, Chang DK, Nones K, et al: Genomic analyses identify molecular subtypes of pancreatic cancer. *Nature* 531:47–52, 2016
- 15.** Raphael BJ, Hruban RH, Aguirre AJ, et al: Integrated Genomic Characterization of Pancreatic Ductal Adenocarcinoma. *Cancer Cell* 32:185–203.e13, 2017
- 16.** Sandhu V, Wedge DC, Bowitz Lothe IM, et al: The Genomic Landscape of Pancreatic and Periampullary Adenocarcinoma. *Cancer Res* 76:5092–5102, 2016
- 17.** Van den Broeck A, Vankelecom H, Van Delm W, et al: Human pancreatic cancer contains a side population expressing cancer stem cell-associated and prognostic genes. *PLoS One* 8:e73968, 2013
- 18.** Donahue TR, Tran LM, Hill R, et al: Integrative survival-based molecular profiling of human pancreatic cancer. *Clin Cancer Res* 18:1352–1363, 2012
- 19.** Sergeant G, van Eijnsden R, Roskams T, et al: Pancreatic cancer circulating tumour cells express a cell motility gene signature that predicts survival after surgery. *BMC Cancer* 12:527, 2012
- 20.** Newhook TE, Blais EM, Lindberg JM, et al: A thirteen-gene expression signature predicts survival of patients with pancreatic cancer and identifies new genes of interest. *PLoS One* 9:e105631, 2014
- 21.** Stratford JK, Bentrem DJ, Anderson JM, et al: A six-gene signature predicts survival of patients with localized pancreatic ductal adenocarcinoma. *PLoS Med*



7:e1000307, 2010

- 22.** Birnbaum DJ, Finetti P, Lopresti A, et al: A 25-gene classifier predicts overall survival in resectable pancreatic cancer. *BMC Med* 15:170, 2017
- 23.** Chen D-T, Davis-Yadley AH, Huang P-Y, et al: Prognostic Fifteen-Gene Signature for Early Stage Pancreatic Ductal Adenocarcinoma. *PLoS One* 10:e0133562, 2015
- 24.** Haider S, Wang J, Nagano A, et al: A multi-gene signature predicts outcome in patients with pancreatic ductal adenocarcinoma. *Genome Med* 6:105, 2014
- 25.** Patil P, Bachant-Winner P-O, Haibe-Kains B, et al: Test set bias affects reproducibility of gene signatures. *Bioinformatics* 31:2318–2323, 2015
- 26.** Eddy JA, Sung J, Geman D, et al: Relative expression analysis for molecular cancer diagnosis and prognosis. *Technol Cancer Res Treat* 9:149–159, 2010
- 27.** Tan AC, Naiman DQ, Xu L, et al: Simple decision rules for classifying human cancers from gene expression profiles. *Bioinformatics* 21:3896–3904, 2005
- 28.** Afsari B, Fertig EJ, Geman D, et al: switchBox: an R package for k-Top Scoring Pairs classifier development. *Bioinformatics* 31:273–274, 2015
- 29.** Harrell FE Jr, Califf RM, Pryor DB, et al: Evaluating the yield of medical tests. *JAMA* 247:2543–2546, 1982
- 30.** Cochran WG: The Combination of Estimates from Different Experiments.

Biometrics 10:101–129, 1954

**31.** Schröder MS, Culhane AC, Quackenbush J, et al: survcomp: an R/Bioconductor package for performance assessment and comparison of survival models. *Bioinformatics* 27:3206–3208, 2011

**32.** Haibe-Kains B, Desmedt C, Sotiriou C, et al: A comparative study of survival models for breast cancer prognostication based on microarray data: does a single gene beat them all? *Bioinformatics* 24:2200–2208, 2008

**33.** Royston P, Sauerbrei W: A new measure of prognostic separation in survival data. *Stat Med* 23:723–748, 2004

**34.** Harrell FE Jr, Lee KL, Mark DB: Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors. *Stat Med* 15:361–387, 1996

**35.** Kassambara A, Kosinski M, Biecek P: survminer: Drawing Survival Curves using 'ggplot2'. R package version 0.3.1, 2017

**36.** Våremo L, Nielsen J, Nookaew I: Enriching the gene set analysis of genome-wide data by incorporating directionality of gene expression and combining statistical hypotheses and methods. *Nucleic Acids Res* 41:4378–4391, 2013

**37.** Liberzon A, Subramanian A, Pinchback R, et al: Molecular signatures database (MSigDB) 3.0. *Bioinformatics* 27:1739–1740, 2011

**38.** Liberzon A, Birger C, Thorvaldsdóttir H, et al: The Molecular Signatures

Database (MSigDB) hallmark gene set collection. *Cell Syst* 1:417–425, 2015

**39.** Benjamini Y, Hochberg Y: Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J R Stat Soc Series B Stat Methodol* 57:289–300, 1995

**40.** Gendoo DMA, Ratanasirigulchai N, Schröder MS, et al: Genefu: an R/Bioconductor package for computation of gene expression-based signatures in breast cancer. *Bioinformatics* 32:1097–1099, 2016

**41.** Sandve GK, Nekrutenko A, Taylor J, et al: Ten simple rules for reproducible computational research. *PLoS Comput Biol* 9:e1003285, 2013

**42.** Gentleman R: Reproducible research: a bioinformatics case study. *Stat Appl Genet Mol Biol* 4:Article2, 2005

**43.** Stroup DF, Berlin JA, Morton SC, et al: Meta-analysis of Observational Studies in Epidemiology: A Proposal for Reporting. *JAMA* 283:2008–2012, 2000

**44.** Notta F, Chan-Seng-Yue M, Lemire M, et al: A renewed model of pancreatic cancer evolution based on genomic rearrangement patterns. *Nature* 538:378–382, 2016

**45.** Kirby MK, Ramaker RC, Gertz J, et al: RNA sequencing of pancreatic adenocarcinoma tumors yields novel expression patterns associated with long-term survival and reveals a role for ANGPTL4. *Mol Oncol* 10:1169–1182, 2016

**46.** Nones K, Waddell N, Song S, et al: Genome-wide DNA methylation patterns in

pancreatic ductal adenocarcinoma reveal epigenetic deregulation of SLIT-ROBO, ITGA2 and MET signaling. *Int J Cancer* 135:1110–1118, 2014

**47.** Sandhu V, Bowitz Lothe IM, Labori KJ, et al: Molecular signatures of mRNAs and miRNAs as prognostic biomarkers in pancreatobiliary and intestinal types of periampullary adenocarcinomas. *Mol Oncol* 9:758–771, 2015

**48.** Zhang G, Schetter A, He P, et al: DPEP1 inhibits tumor cell invasiveness, enhances chemosensitivity and predicts clinical outcome in pancreatic ductal adenocarcinoma. *PLoS One* 7:e31507, 2012

**49.** Winter C, Kristiansen G, Kersting S, et al: Google goes cancer: improving outcome prediction for cancer patients by network-based ranking of marker genes. *PLoS Comput Biol* 8:e1002511, 2012

**50.** Zheng X, Carstens JL, Kim J, et al: Epithelial-to-mesenchymal transition is dispensable for metastasis but induces chemoresistance in pancreatic cancer. *Nature* 527:525–530, 2015

**51.** Gaianigo N, Melisi D, Carbone C: EMT and Treatment Resistance in Pancreatic Cancer [Internet]. *Cancers* 9, 2017 Available from: <http://dx.doi.org/10.3390/cancers9090122>

**52.** FUKUSHIMA, N: Ductal adenocarcinoma variants and mixed neoplasm of the pancreas. *WHO Classification of Tumours of the Digestive System* 292–295, 2010



## FIGURES

**Figure 1.** Pipeline showing the approach used for building the Pancreatic Cancer Overall Survival Predictor (PCOSP).

**Figure 2:** Flowchart showing the inclusion criteria for the pancreatic adenocarcinoma samples.

**Figure 3. Predictive value of PCOSP for early death and overall survival. (A)** Area under the ROC curves for all the cohorts and the meta estimates for sequencing cohorts, array-based cohorts and for both the platforms combined. Forestplot reporting **(B)** the concordance indices (C-index) and **(C)** the hazard ratio (HR) for all the cohorts and the meta estimates for sequencing cohorts (orange), array-based cohorts (blue) and for both the platforms combined (grey).

**Figure 4. Kaplan Meier survival curves** for **(A)** PCSI **(B)** TCGA **(C )** Kirby **(D)** ICGC-array **(E)** UNC **(F)** Chen **(G)** OUH **(H)** Zhang **(I)** Winter and **(J)** Collisson. The overall survival difference between high and low risk group is 13 and 23 months respectively.

**Figure 5. Comparison of the prognostic value of the clinicopathological model and PCOSP. (A)** Barplot reporting the AUROCs for the clinical model and the PCOSP model. (Forestplot reporting the the **(B)** concordance index (C-index) and **(C)** Hazard ratio (HR) of validation cohorts computed using PCOSP, and clinicopathological model.

**Figure 6. Comparison of existing classifiers with PCOSP.** The forestplot reports the meta-estimate of **(A)** concordance indices (C-index) and **(B)** hazard ratio (HR) for PCOSP and existing classifiers.

## **SUPPLEMENTARY FIGURES**

**Supplementary Figure S1: Density plot showing the distribution of balanced accuracy for random models.** Distribution of meta-estimates of 1000 models generated using **(A)** random reshuffling of labels and **(B)** random assignment of genes to TSP models. The meta-estimates were independently calculated for all the cohorts combined, sequencing cohorts and array-based cohort. The pink, green, blue dashed lines represent meta-estimate of AUROC from PCOSP model for overall, sequencing and array-based cohorts respectively.

**Supplementary Figure S2: Forestplot of (A) concordance index (C-index) and (B) hazard ratio (HR) for all the cohorts divided based on the molecular subtypes.** The grey, green and pink color in the forestplot depicts meta-estimate of C-index for overall cohort, the basal subtype and the classical subtype of the cohorts, respectively.

**Supplementary Figure S3:** The scatterplot shows the meta-estimate of AUROC (orange) and total number of unique genes (blue) in the PCOSP model at different

balanced accuracy thresholds. The threshold used in the PCOSP is marked as dashed line at 0.6.

## **SUPPLEMENTARY TABLES**

**Supplementary Table S1:** The table shows the datasets used in the project for meta-analysis.

**Supplementary Table S2:** The table shows the clinicopathological information of the validation cohorts used in the analysis.

**Supplementary Table S3:** Univariate and multivariate regression analysis from (A) logistic regression model to predict early death (death >1 yr or not), and (B). the Cox regression model using clinicopathological features, molecular subtypes and PCOSP model probabilities for validation cohorts.

**Supplementary table S4:** The table shows the pathways overrepresented in the PCOSP model genes using (A) hallmark gene sets, (B) canonical pathways, (C) GO-molecular function term (and (D) GO cellular component terms from MSigDB.